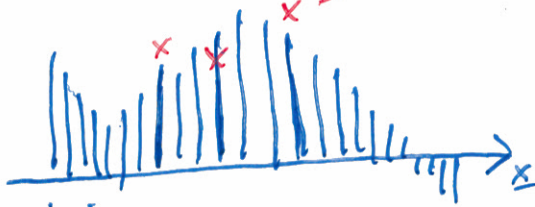


# Gaussian Processes



Noisy observations of a subset values.

- Can infer any subset of function values given these.

Gaussian:

Prior directly on ~~the~~ the function values at all possible inputs

$$\text{Marginal } p(\underline{f}) = N(\underline{f}; \underline{0}, \mathbf{K})$$

↑  
vector of values at locations  $\{\underline{x}^{(i)}\}$

$$\begin{aligned} \uparrow \\ K_{ij} &= k(\underline{x}^{(i)}, \underline{x}^{(j)}) \\ &= \text{cov}[f_i, f_j] \end{aligned}$$

# GrPs are Bayesian Linear Regression

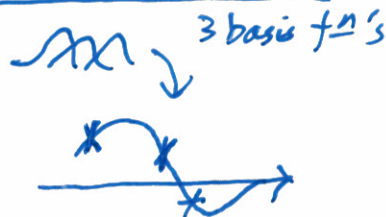
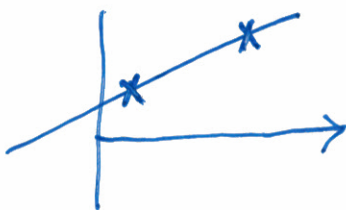
Model

$$f_i = f(\underline{x}^{(i)}) = \underline{w}^T \underline{x}^{(i)} + b$$

Prior:  $\underline{w} \sim N(0, \sigma_w^2 I), b \sim N(0, \sigma_b^2)$

$\uparrow$   
( $\sigma_{\text{prior}}, \frac{1}{\alpha}$ )

$$\begin{aligned} \text{cov}(f_i, f_j) &= \mathbb{E}[f_i f_j] - \underbrace{\mathbb{E}[f_i]}_0 \underbrace{\mathbb{E}[f_j]}_0 \\ &= \mathbb{E}\left[\underbrace{\underline{x}^{(i)T} \underline{w}}_{\underline{w}^T \underline{x}^{(i)}} + b\right) \left(\underline{w}^T \underline{x}^{(j)} + b\right)\right] \\ &= \mathbb{E}\left[\underline{x}^{(i)T} \underline{w} \underline{w}^T \underline{x}^{(j)} + b^2 + \dots\right] \\ &= \underline{x}^{(i)T} \underbrace{\mathbb{E}[\underline{w} \underline{w}^T]}_{\text{cov}(\underline{w})} \underline{x}^{(j)} + \underbrace{\mathbb{E}[b^2]}_0 + \dots \\ &= \sigma_w^2 \underline{x}^{(i)T} \underline{x}^{(j)} + \sigma_b^2 = k(\underline{x}^{(i)}, \underline{x}^{(j)}) \end{aligned}$$



## Basis Functions

$$k(\underline{x}^{(i)}, \underline{x}^{(j)}) = \sigma_w^2 \underbrace{\phi(\underline{x}^{(i)})^T \phi(\underline{x}^{(j)})}_{\text{inner product of features}} + \sigma_b^2$$

We only need inner products of features.

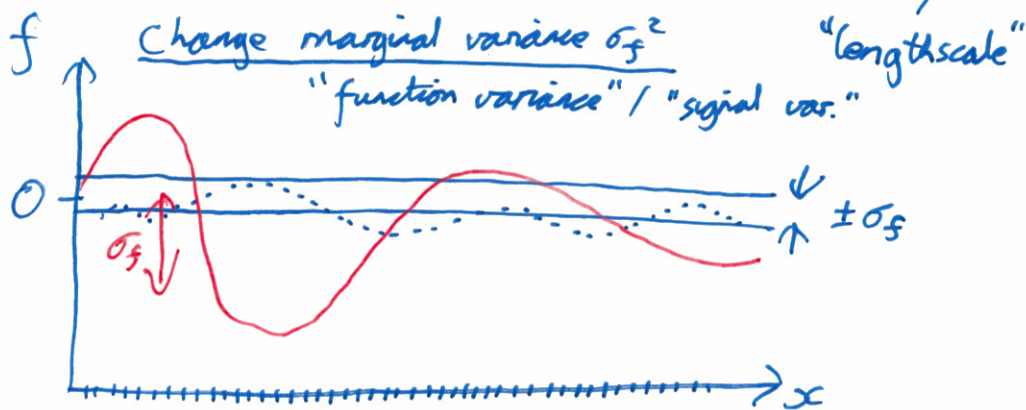
## "kernel trick"

- Rewrite algorithm so it only needs inner products of features.
- We then use very large / infinite set of basis functions.
- Replace that inner product with analytic expression we can compute.

It can be shown that....

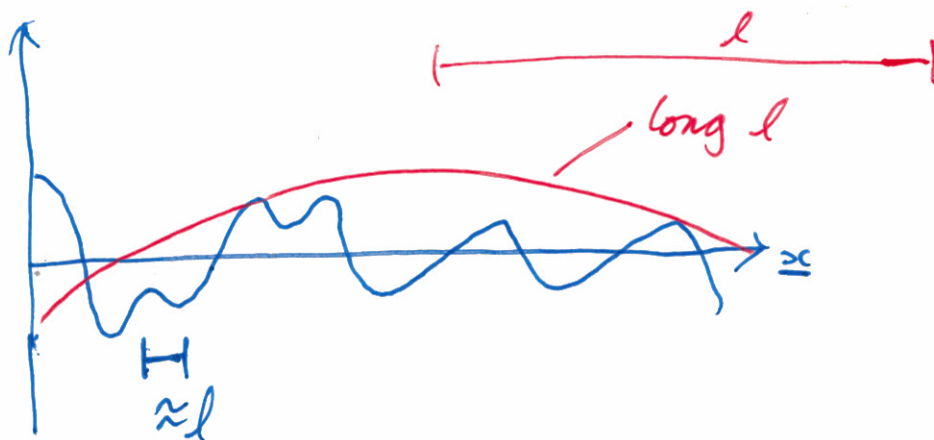
If we put RBFs everywhere we can derive a kernel:

$$k(\underline{x}^{(i)}, \underline{x}^{(j)}) = \sigma_f^2 \exp\left(-\frac{1}{2} \sum_d (x_d^{(i)} - x_d^{(j)})^2 / l_d^2\right)$$



$$k(\underline{x}^{(i)}, \underline{x}^{(i)}) = \sigma_f^2 \quad \text{Diagonal of cov. matrix}$$

## Change lengthscale



Pick parameters by (marginal) likelihood.

$$P(y | X, \theta = \{\sigma_y^2, l, \sigma_n^2, \dots\})$$

$$= N(y; 0, K + \sigma_n^2 I).$$

↑

$O(N^2)$  RAM

Inverting / Factoring

costs  $O(N^3)$ ...